



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Leveraging hierarchical memories for micro-core architectures

Citation for published version:

Brown, N & Jamieson, M 2018, 'Leveraging hierarchical memories for micro-core architectures', 5th International Conference on Exascale Applications and Software, Edinburgh, United Kingdom, 17/04/18 - 19/04/18.

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Publisher's PDF, also known as Version of record

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



Leveraging hierarchical memories for micro-core architectures

Nick Brown,¹ Maurice Jamieson¹

¹*EPCC, The University of Edinburgh, Scotland, n.brown@epcc.ed.ac.uk, maurice.jamieson@ed.ac.uk*

1. Background

Micro-core architectures combine many simple, low power and low on-chip memory cores onto a single processor package. The low power nature of these architectures means that there is potential for their use in future HPC and embedded systems, and their low cost makes them ideal for education and prototyping. However there is a high barrier to entry in programming due to the considerable complexity and immaturity of supporting tools.

ePython is a Python virtual machine we have developed for the 16-core Epiphany III micro-core architecture which fits in the 32Kb per core memory. In combination we developed an abstraction that supports offloading functions in existing Python codes, running on a host CPU, seamlessly to the micro-cores. In [1] we introduced this abstraction and motivated it with a machine learning code for detecting lung cancer in 3D CT scans where kernels for model training and inference ran in parallel on the micro-cores. However the small amount of core memory severely limited the physical size of the images, which had to be interpolated to fit.

2. Hierarchical memories

In addition to the small on-core memory, there is typically much larger, slower external memory as illustrated for the Epiphany by figure 1. In order to take full advantage of these architectures one must leverage these hierarchies of memory, but a key question is how best to achieve this whilst maintaining good performance.

In this work we have addressed this challenge by splitting the memory abstraction into three choices:

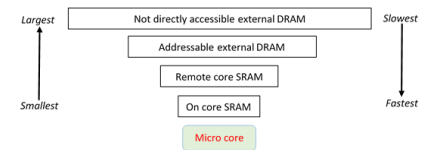


Figure 1. Typical Epiphany memory hierarchy

- A mirroring of memory where copies of external memory also exist on the micro-cores. A manually copying of data to and from the different memory levels is required.
- Memory can be exposed from a specific level in the hierarchy to the micro-cores without an explicit copy being allocated. Reads and writes directly access this external memory, these accesses being blocking or non-blocking (using the DMA engines.) Abstractions around the non-blocking approach enables the programmer to leverage patterns such as double buffering and data streaming to overlap compute and memory access for performance.
- By default memory belongs to the hierarchical level where it is first declared. It is possible to override this via memory kinds [2]. In our approach these are Python objects that follow a standard interface and sit outside of the core ePython implementation. They define the behaviour of memory access at the level of hierarchy they represent.

Based upon this work we are now able to run the machine learning code of [1] with the full sized images. The programmer is able to experiment with choices around memory placement and access patterns without having to worry about the low level complexities of data movement.

References

- [1] Brown, N. Offloading Python kernels to micro-core architectures, Poster Supercomputing 2017, Denver US
- [2] Cantalupo, C. et al. Memkind: An Extensible Heap Memory Manager for Heterogeneous Memory Platforms and Mixed Memory Policies. No. SAND2015-1862C. SNL-NM, Albuquerque, NM (United States), 2015.